



MAX-PLANCK-GESELLSCHAFT

Sparse Gaussian Process Toolbox

Lehel Csató

Max Planck Institute for Biological Cybernetics, Tübingen, Germany
Neural Computing Research Group, Aston University, Birmingham



BIOLOGISCHE KYBERNETIK

Toolbox Overview

Uses:

- Gaussian process (\mathcal{GP}) latent model
- Sequential approximation to the posterior
- Sparsification of the resulting process
- *MATLAB* programming language
- *NETLAB* toolbox

Provides:

- GUI demos for teaching \mathcal{GP} s
- Variety of error or likelihood functions
- Bayesian hyperparameter selection

Freely available from:

- <http://www.tuebingen.mpg.de/~csatol>
- <http://www.ncrg.aston.ac.uk/Projects/SSGP>

Gaussian Process Inference

Gaussian process (\mathcal{GP}) models are *probabilistic* kernel methods. \mathcal{GP} s specify priors over a function space. Any finite sample from the random function has joint Gaussian distribution with covariance given by a kernel function. The *prior* is thus

$$p_0(\mathbf{f}) = \frac{1}{(2\pi)^{N/2} |\mathbf{K}_N|^{1/2}} \exp\left(-\frac{1}{2} \mathbf{f}^T \mathbf{K}_N^{-1} \mathbf{f}\right) \quad (1)$$

where

$$\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^T \text{ are samples from the function}$$

$$\mathbf{K}_N = \{K_0(\mathbf{x}_i, \mathbf{x}_j; \theta)\}_{i,j=1}^N \text{ is the sample covariance matrix}$$

where $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ are the inputs and θ is the set of parameters of K_0 . For common \mathcal{GP} models the likelihoods factorise:

$$P(\mathcal{D}) = \prod_{n=1}^N t(y_n | \mathbf{f}) \doteq \prod_{n=1}^N t_n(\mathbf{f})$$

The posterior *process* is

$$p_{\text{post}}(\mathbf{f}) = \frac{1}{Z} p_0(\mathbf{f}) P(\mathcal{D} | \mathbf{f})$$

For non-Gaussian likelihoods the posterior is not a \mathcal{GP} , thus it is analytically intractable.

- There are **problems of:** (solved by) *variational appr.*
- tractability for non-Gaussian likelihoods:
 - representation for large datasets: *sparse appr.*

Approximations to the Posterior

Representation of the posterior moments

Using \mathcal{GP} priors, the moments of the posterior process are:

$$\langle f(\mathbf{x}) \rangle_{\text{post}} = \langle f(\mathbf{x}) \rangle_0 + \sum_{i=1}^N K_0(\mathbf{x}, \mathbf{x}_i) \alpha(i) \quad (2)$$

$$K(\mathbf{x}, \mathbf{x}')_{\text{post}} = K_0(\mathbf{x}, \mathbf{x}') + \sum_{i,j=1}^N K_0(\mathbf{x}, \mathbf{x}_i) C(i,j) K_0(\mathbf{x}_j, \mathbf{x}')$$

where $\alpha(i)$ and $C(i,j)$ are parameters driven by the likelihood function. The representation suggests an approximation: the posterior is approximated by the closest \mathcal{GP} in a KL-sense. The process approximation is reduced to finding the parameters $\alpha = [\alpha(1), \dots, \alpha(N)]^T$ and $C = \{C(i,j)\}_{i,j=1}^N$.

Approximation to the posterior process

The approximations are sequential, including a single case at each step. We use the *expectation-propagation (or TAP) algorithm* where the posterior is constructed from local Gaussian approximation $\hat{t}_i(\mathbf{f})$ to each factor in the likelihood [Minka 2000; Opper and Winther 2001].

We define the approximating process:

$$p_{\text{post}}(\mathbf{f}) \approx \hat{p}(\mathbf{f}) \propto p_0(\mathbf{f}) \prod_{n=1}^N \hat{t}_n(\mathbf{f})$$

and compute it using the algorithm:

- Initialise $t_n(\mathbf{f}) = 1 \quad i = 1, \dots, N$.
- For each input n compute: $p_0^n(\mathbf{f}) \propto \frac{\hat{p}(\mathbf{f})}{t_n(\mathbf{f})}$
- Approximate the *local* posterior and substitute back $\hat{t}_n(\mathbf{f})$ based on:

$$\frac{1}{Z_n} p_0^n(\mathbf{f}) t_n(\mathbf{f}) \approx p_{\text{post}}^n(\mathbf{f}) = \frac{1}{Z_n} p_0^n(\mathbf{f}) \hat{t}_n(\mathbf{f}) \Rightarrow t_n(\mathbf{f}) \approx \frac{Z_n}{Z_n} \hat{t}_n(\mathbf{f}) \quad (3)$$

At the equilibrium point we have an approximation to the *marginal likelihood* (or evidence) as:

$$Z = \int d\mathbf{f} p_0(\mathbf{f}) \prod_{n=1}^N t_n(\mathbf{f}) \approx \prod_{n=1}^N \int d\mathbf{f} p_0(\mathbf{f}) \hat{t}_n(\mathbf{f}) \quad (4)$$

Sparsification

If $t_n(\mathbf{f})$ depends only on f_n , then \mathbf{f} reduces to f_n . The random variable f_n can further be eliminated by

$$f_n \rightarrow \hat{f}_n = \pi_n \mathbf{f}_B \quad \text{, where } B \text{ is a predefined set of inputs}$$

(B can be from the training/test data) The approximated posterior \mathcal{GP} has the mean function and covariance kernel defined as:

$$\langle f(\mathbf{x}) \rangle_{\text{post}} = \langle f(\mathbf{x}) \rangle_0 + \sum_{i \in B} K_0(\mathbf{x}, \mathbf{x}_i) \hat{\alpha}(i)$$

$$K(\mathbf{x}, \mathbf{x}')_{\text{post}} = K_0(\mathbf{x}, \mathbf{x}') + \sum_{i,j \in B} K_0(\mathbf{x}, \mathbf{x}_i) C(i,j) K_0(\mathbf{x}_j, \mathbf{x}')$$

Important: user control over the size of B. \Rightarrow possible to use large datasets.

Matlab Implementation

Uses a *matlab structure net*. The initialisation of the structure with the default values for the fields is done using:

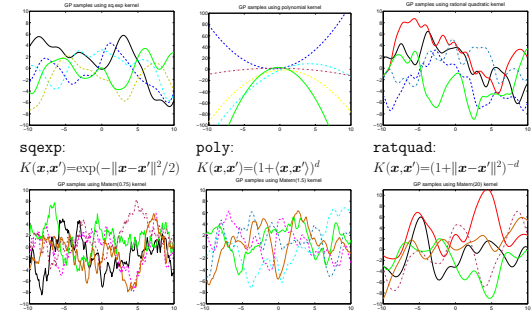
```
net = ogp(odim, odim, covarfn, covpar);
```

where o_{dim} is the dimension of inputs, o_{dim} is the dimension of outputs – one can define pairs of \mathcal{GP} s for more than a single latent variable in the likelihood. Associated to each input dimension there is an ARD [MacKay 1992] parameter: `log(v)=net.inweights` and the kernel functions depend on the weighted product:

$$(\mathbf{x}, \mathbf{x}') = \sum_i \gamma_i x_i x'_i \quad \text{which is used in the kernels below}$$

Kernel functions for the prior process

`covarfn`, `covpar` define the covariance function for the \mathcal{GP} . The kernel hyperparameters are $\theta = \text{covpar}$. The covariance functions are implemented:



user - extensibility of the toolbox with user-defined covariance function (example: **Matérn kernel**) [Stein 1999]:

$$K(\mathbf{x}, \mathbf{x}') = \frac{A}{\Gamma(\nu) 2^{(\nu-1)}} (\sqrt{2\nu d} K_\nu(\sqrt{2\nu d}))$$

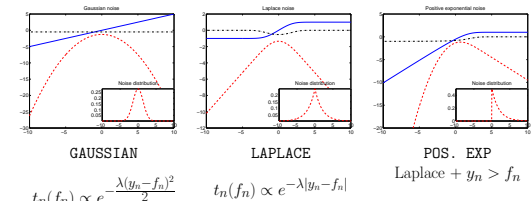
where K_ν is the modified Bessel function of the second kind. Allows transition from rough covariances to squared exponentials by $\nu \rightarrow \infty$.

Likelihood models

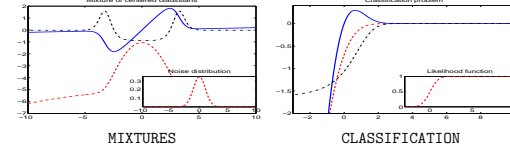
The toolbox requires a function which returns the local update coefficients from eq. (3). It is specified by `net = ogpinit(net, @likfn, likpar)`

For Gaussian noise model (`c_reg_gauss`) no approximation is needed.

Implemented likelihood models



Obs: Nonstandard likelihood models are difficult to deal with using conventional kernel methods.



Inference and prediction

Iterating the following two-steps (EM algorithm):

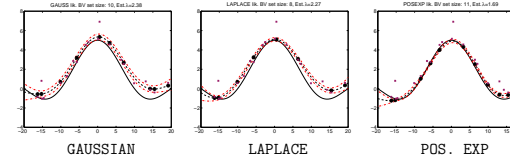
- Given the set of kernel- and likelihood parameters, find \mathcal{GP}_{opt} using the *EP* algorithm.
- Fixing the \mathcal{GP} , optimising the *evidence* from eq. (4) with respect to hyperparameters. For any kernel parameter θ

$$\frac{\partial \ln Z}{\partial \theta} = \text{tr} \left[\frac{\partial \ln Z}{\partial \mathbf{K}_B} \frac{\partial \mathbf{K}_B}{\partial \theta} \right]$$

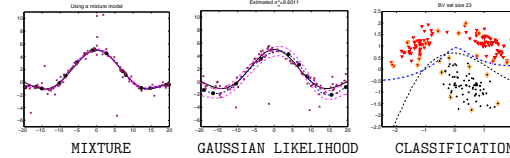
```
net = ogphyplearn(net, its)
```

Examples

Noise generation: **Posexp** with $\lambda = 1.66$.



Using mixtures to tackle outliers:



Acknowledgements

Comments and discussions with Dan Cornford (NCRG), Ian Nabney (NCRG), Carl Rasmussen (MPI), codes from Ian Nabney (NCRG) and Anton Schwaighofer (TU Graz) are acknowledged.

References

- MacKay, D. J. C. (1992). Bayesian interpolation. *Neural Computation* **4**, 415–447.
- Minka, T. P. (2000). *Expectation Propagation for Approximate Bayesian Inference*. Ph. D. thesis, Dep. of El. Eng. & Comp. Sci., MIT, [vismod.www.media.mit.edu/~tppminka](http://www.media.mit.edu/~tppminka).
- Opper, M. and O. Winther (2001). Adaptive and self-averaging TAP mean field theory for probabilistic modeling. *Physical Review E* **64**(056131), 1–14.
- Stein, M. L. (Ed.) (1999). *Interpolation of Spatial Data: Some theory for kriging*. New York: Springer.